

PAPER • OPEN ACCESS

Prediction of chemical reaction yields using deep learning

To cite this article: Philippe Schwaller *et al* 2021 *Mach. Learn.: Sci. Technol.* **2** 015016

View the [article online](#) for updates and enhancements.

You may also like





- [A new formulation of gradient boosting](#)
Alex Wozniakowski, Jayne Thompson, Mile Gu et al.
- [Data-driven discovery of Koopman eigenfunctions for control](#)
Eurika Kaiser, J Nathan Kutz and Steven L Brunton
- [How isotropic kernels perform on simple invariants](#)
Jonas Paccolat, Stefano Spigler and Matthieu Wyart



PAPER

Prediction of chemical reaction yields using deep learning

OPEN ACCESS

Philippe Schwaller^{1,2} , Alain C Vaucher¹ , Teodoro Laino¹  and Jean-Louis Reymond² RECEIVED
3 August 2020¹ IBM Research—Europe, Säumerstrasse 4, 8803 Rüschlikon, SwitzerlandREVISED
8 October 2020² Department of Chemistry and Biochemistry, University of Bern, Freiestrasse 3, 3012 Bern, SwitzerlandACCEPTED FOR PUBLICATION
5 November 2020E-mail: pshs@zurich.ibm.comPUBLISHED
31 March 2021**Keywords:** chemical reactions, yield prediction, deep learning, transformerSupplementary material for this article is available [online](#)

Original Content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](#).

Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

**Abstract**

Artificial intelligence is driving one of the most important revolutions in organic chemistry. Multiple platforms, including tools for reaction prediction and synthesis planning based on machine learning, have successfully become part of the organic chemists' daily laboratory, assisting in domain-specific synthetic problems. Unlike reaction prediction and retrosynthetic models, the prediction of reaction yields has received less attention in spite of the enormous potential of accurately predicting reaction conversion rates. Reaction yields models, describing the percentage of the reactants converted to the desired products, could guide chemists and help them select high-yielding reactions and score synthesis routes, reducing the number of attempts. So far, yield predictions have been predominantly performed for high-throughput experiments using a categorical (one-hot) encoding of reactants, concatenated molecular fingerprints, or computed chemical descriptors. Here, we extend the application of natural language processing architectures to predict reaction properties given a text-based representation of the reaction, using an encoder transformer model combined with a regression layer. We demonstrate outstanding prediction performance on two high-throughput experiment reactions sets. An analysis of the yields reported in the open-source USPTO data set shows that their distribution differs depending on the mass scale, limiting the data set applicability in reaction yields predictions.

1. Introduction

Chemical reactions in organic chemistry are described by writing the structural formula of reactants and products separated by an arrow, representing the chemical transformation by specifying how the atoms rearrange between one or several reactant molecules and one or several product molecules [1]. Economic, logistic, and energetic considerations drive chemists to prefer chemical transformations capable of converting all reactant molecules into products with the highest yield possible. However, side-reactions, degradation of reactants, reagents or products in the course of the reaction, equilibrium processes with incomplete conversion to a product, or simply by product isolation and purification undermine the quantitative conversion of reactants into products, rarely reaching optimal performance.

Reaction yields are usually reported as a percentage of the theoretical chemical conversion, i.e. the percentage of the reactant molecules successfully converted to the desired product compared to the theoretical value. It is not uncommon for chemists to synthesise a molecule in a dozen or more reaction steps. Hence, low-yield reactions may have a disastrous effect on the overall route yield because of the individual steps' multiplicative effect. Therefore, it is not surprising that designing new reactions with yields higher than existing ones attracts much effort in organic chemistry research.

In practice, specific chemical reaction classes are characterised by lower or higher yields, with the actual value depending on the reaction conditions (temperature, concentrations, etc) and on the specific substrates.

Estimating the reaction yield can be a game-changing asset for synthesis planning. It provides chemists with the ability to evaluate the overall yield of complex reaction paths, addressing possible shortcomings well ahead of investing hours and materials in wet-lab experiments. Computational models predicting reaction

yields could support synthetic chemists in choosing an appropriate synthesis route among many predicted by data-driven algorithms. Moreover, reaction yields prediction models could also be employed as scoring functions in computer-assisted retrosynthesis route planning tools [2–5], to complement forward prediction models [4, 6] and in-scope filters [2].

Most of the existing efforts in constructing models for the prediction of reactivity or of reaction yields focused on a particular reaction class: oxidative dehydrogenations of ethylbenzene with tin oxide catalysts [7], reactions of vanadium selenites [8], Buchwald–Hartwig aminations [9–11], and Suzuki–Miyaura cross-coupling reactions [12–14]. To the best of our knowledge, there has been only one attempt to design a general-purpose prediction model for reactivity and yields, without applicability constraints to a specific reaction class [15]. In this work, the authors design a model predicting whether the reaction yield is above or below a threshold value and conclude that the models and descriptors they consider cannot deliver satisfactory results.

Here, we build on our legacy of treating organic chemistry as a language to introduce a new model that predicts reaction yields starting from reaction SMILES [16]. More specifically, we fine-tune the rxnfp models by Schwaller *et al* [17] based on a bidirectional encoder representations from transformers (BERT)-encoder [18] by extending it with a regression layer to predict reaction yields. BERT encoders belong to the transformer model family, which has revolutionised natural language processing [18, 19]. These models take sequences of tokens as input to compute contextualised representations of all the input tokens, and can be applied to reactions represented in the SMILES [20] format. In this work, we demonstrate, for the first time, that these natural language architectures are very useful not only when working with language tokens but also in providing descriptors of high quality to predict reaction properties such as reaction yields.

It is possible to train our approach both on data specific to a given reaction class or on data representing different reaction types. Thus, we initially trained the model on two high-throughput experimentation (HTE) data sets. Among the few HTE reaction data sets published in recent years, we selected the data sets for palladium-catalysed Buchwald–Hartwig reactions provided by Ahneman *et al* [9] and for Suzuki–Miyaura coupling reactions provided by Perera *et al* [21]. Finally, we trained our model on patent data available in the USPTO data set [22, 23].

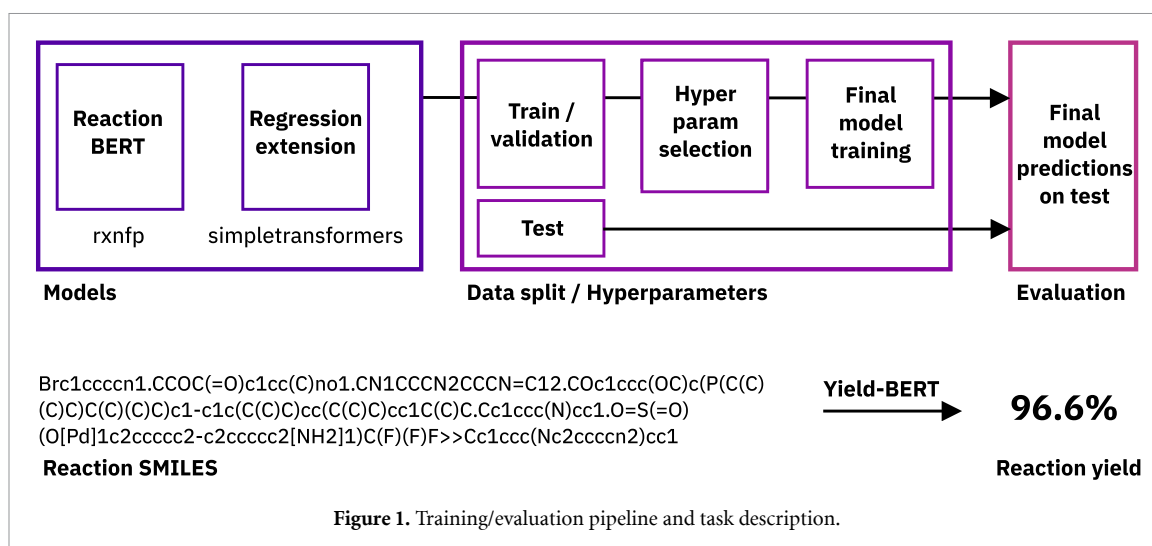
HTE and patent data sets are very different in terms of content and quality. HTE data sets typically cover a very narrow region in the chemical reaction space, with chemical reaction data related to one or a few reaction templates applied to large combinations of selected precursors (reactants, solvents, bases, catalysts, etc). In contrast, patent reactions cover a much wider reaction space. In terms of quality, HTE data sets report reactions represented uniformly and with yields measured using the same analytical equipment, thus providing a consistent and high-quality collection of knowledge. In comparison, the yields from patents were measured by different scientists using different equipment. Incomplete information in the original documents, such as unreported reagents or reaction conditions, and the extensive limitation in text mining technologies makes the entire set of patent reactions quite noisy and sparse. An extensive analysis of the USPTO data set revealed that the experimental conditions and reaction parameters, such as scale of the reaction, concentrations, temperature, pressure, or reaction duration, may have a significant effect on the measured reaction yields. The functional dependency of the yields from the reaction conditions poses additional constraints, as the model presented in this work does not consider those values explicitly in the reaction descriptor. The basic assumption is that every reaction yield reported in the data set is optimised for the reaction parameters.

Our best-performing model reached an R^2 score of 0.956 on a random split of the Buchwald–Hartwig data set, while the highest R^2 score on the smoothed USPTO data was 0.388. These numbers reflect how the intrinsic data set limitations increase the complexity of training a sufficiently good performing model on the patent data, resulting in a more difficult challenge than training a model for the HTE data set.

2. Models and experimental pipeline

We base our models directly on the reaction fingerprint (rxnfp) models by Schwaller *et al* [17]. We use a fixed-size encoder model size, tuning only the hyperparameter for dropout rate and learning rate, thus avoiding often-encountered difficulties of neural networks with numerous hyperparameters. During our experiments, we observed good performances for a wide range of dropout rates (from 0.1 to 0.8) and conclude that the initial learning rate is the most important hyperparameter to tune. Figures S26–S30 show hyperparameter optimisation plots (available online at stacks.iop.org/MLST/2/015016/mmedia). To facilitate the training, our work uses simpletransformers [24], a huggingface transformer [25] and the PyTorch framework [26]. The overall pipeline is shown in figure 1.

To provide an input compatible with the rxnfp model we use the same RDKit [27] reaction canonicalisation and SMILES tokenization [6] as in the rxnfp work [17].



3. High-throughput experiment yield predictions

3.1. Buchwald–Hartwig reactions

Ahneman *et al* [9] performed high-throughput experiments on Pd-catalysed Buchwald–Hartwig C–N cross coupling reactions, measuring the yields for each reaction. For the experiments, they used three 1536-well plates spanning a matrix of 15 aryl and heteroaryl halides, four Buchwald ligands, three bases, and 23 isoxazole additives, resulting in 3955 reactions. As inputs for their models, Ahneman *et al* [9] computed 120 molecular, atomic and vibrational properties with density functional theory using Spartan for every halide, ligand, base and additive combination. The descriptors included HOMO and LUMO energy, dipole moment, electronegativity, electrostatic charge and NMR shifts for atoms shared by the reagents. Compared to reaction SMILES that can vary in length, the input in the work of Ahneman *et al* [9] was a fixed-size vector. They investigated numerous methods, including linear models, k-nearest-neighbours, support vector machines, Bayes generalised linear models, artificial neural networks and random forests. Eventually, they selected their random forest model as the best performing one. The work of Ahneman *et al* [9] was challenged by Chuang and Keiser [10], who pointed out several issues. First, by replacing the computed chemical features with random features of the same length or one-hot encoded vectors Chuang and Keiser got similar performance to the original paper with the chemical features. Therefore, they weakened the original claim that additive features were the most important for the predictions. However, the additive features were on average still estimated to be the most important features by the random forest model when the yields were shuffled [10]. Recently, Sandfort *et al* [11] used a concatenation of multiple molecular fingerprints as an alternative reaction representation to demonstrate superior yield prediction performance compared to one-hot encoding.

Unlike previous work, we directly use the reaction SMILES as input to a BERT-based reaction encoder [17] enriched with a regression layer (Yield-BERT). To investigate the suggested method, we used the same splits as Sandfort *et al* [11]. In contrast, to their work, we used 1/7 of the training set from the first random split as a validation set to select optimal values for the two hyperparameters, namely, learning rate and dropout probability. Once selected, we kept the hyperparameters identical for all the subsequent experiments.

The results are shown in table 1. Using solely a reaction SMILES representation, our method achieves an average R^2 of 0.951 on the random splits and outperforms not only the MFF by Sandfort *et al* [11], but also the chemical descriptors computed with density functional theory (DFT) by Ahneman *et al* [9]. Moreover, for the out-of-sample tests where the isoxazole additives define the splits our method performs on average better than multiple-fingerprint features (MFF) and one-hot descriptors and comparable to the chemical descriptors. As in the work of Sandfort *et al* [11], the test 3 split resulted in the worst model performance. For the rest of the out-of-sample, our method performs better than the others. We also reduced the training set to 5% (197 reactions), 10% (395 reactions) and 20% (791 reactions) and observed that the model learned to reasonably predict yields despite the significantly smaller training set. Detailed results on the different Buchwald–Hartwig test sets are shown in figures S1–S14.

3.2. Suzuki–Miyaura reactions

Perera *et al* [21] used HTE technologies on the class of Suzuki–Miyaura reactions. They considered 15 pairs of electrophiles and nucleophiles, each leading to a different product. For each pair, they varied the ligands

Table 1. Comparing methods on the Buchwald–Hartwig data set. All results shown in this table used the rxnfp pretrained model as base encoder.

R^2	DFT [9]	One-hot [10, 11]	MFF [11]	Yield-BERT
Rand 70/30	0.92	0.89	0.927 ± 0.007	0.951 ± 0.005
Rand 50/50	0.9			0.92 ± 0.01
Rand 30/70	0.85			0.88 ± 0.01
Rand 20/80	0.81			0.86 ± 0.01
Rand 10/90	0.77			0.79 ± 0.02
Rand 5/95	0.68			0.61 ± 0.04
Rand 2.5/97.5	0.59			0.45 ± 0.05
Test 1	0.8	0.69	0.85	0.84 ± 0.01
Test 2	0.77	0.67	0.71	0.84 ± 0.03
Test 3	0.64	0.49	0.64	0.75 ± 0.04
Test 4	0.54	0.49	0.18	0.49 ± 0.05
Avg. 1–4	0.69	0.59	0.60	0.73

Note: The best values are shown in bold.

Table 2. Summary of the average R^2 scores on the Suzuki–Miyaura reactions data set using a Yield-BERT with different base encoders. We used 10 different random folds (70/30).

Base encoder rxnfp [17]	Pretrained	Pretrained	ft	ft
Hyperparameters	Same as 3.1	Tuned	Same as 3.1	Tuned
Random 70/30	0.79 ± 0.01	0.79 ± 0.02	0.81 ± 0.02	0.81 ± 0.01

Note: The best values are shown in bold.

(12 in total), bases (8), and solvents (4), resulting in a total of 5760 measured yields. The same data set was also investigated in the work of Granda *et al* [12].

Here, we first trained our yield prediction models with the same hyperparameters as for the Buchwald–Hartwig reaction experiment above, achieving an R^2 score of 0.79 ± 0.01 . Second, we tuned the dropout probability and learning rate, similarly to the previous experiment, using a split of the training set of the first random split. The resulting hyperparameters were then used for all the splits. The hyperparameter tuning did not lead to better performance compared to the parameters used for the Buchwald–Hartwig reactions. This shows that the models have a stable performance for a wide range of parameters and that they are transferable from one data set to another related data set.

We also compared two different base encoder models that are available from the rxnfp library [17], namely the BERT model pretrained with a masked language modelling task, and the BERT model subsequently fine-tuned on a reaction class prediction task. The results are displayed in table 2. In contrast to the Buchwald–Hartwig data set, where no difference between the two base encoders was observed, the ft model achieves an R^2 score of 0.81 ± 0.01 , outperforming the pretrained base encoder on the Suzuki–Miyaura reactions. Detailed results on the different Suzuki–Miyaura test sets are shown in figures S15–S24.

3.3. Discovery of high-yield reactions with reduced training sets

Granda *et al* [12] proposed training on a random (10%) portion of the original data set to evaluate the rest of the reactions with the purpose of selecting the next reactions to test. Similarly, we trained our models on different fractions of the training set and used them to evaluate the yields of the remaining reactions. The aim here is to evaluate how good the models are at selecting high-yield reactions after having seen a small fraction of randomly chosen reactions.

As can be seen from figure 2, training on only 5% of the reactions already enables a chemist to select some of the highest yielding reactions for the next round of the experiments. With a training set of 10% the yields of the selected reactions are close to the best possible selection marked with ‘ideal’ in the figure. For the Buchwald–Hartwig reaction, using a model trained on 10% of the data set, the 10 reactions from the remaining unseen data set predicted to have the highest yields, have an average yield of $90 \pm 6\%$, compared to the ideal selection of $98.7 \pm 0.9\%$. In contrast, a random selection of 10 reactions would have led to yields of $34 \pm 27\%$. The selection works similarly for the Suzuki–Miyaura reactions.

We performed a purely greedy selection, as we aimed to find highest yielding reactions after one training round. A wider chemical reaction space exploration with a reaction selection using more elaborate uncertainty estimates and an active learning strategy was investigated by Eyke *et al* [14].

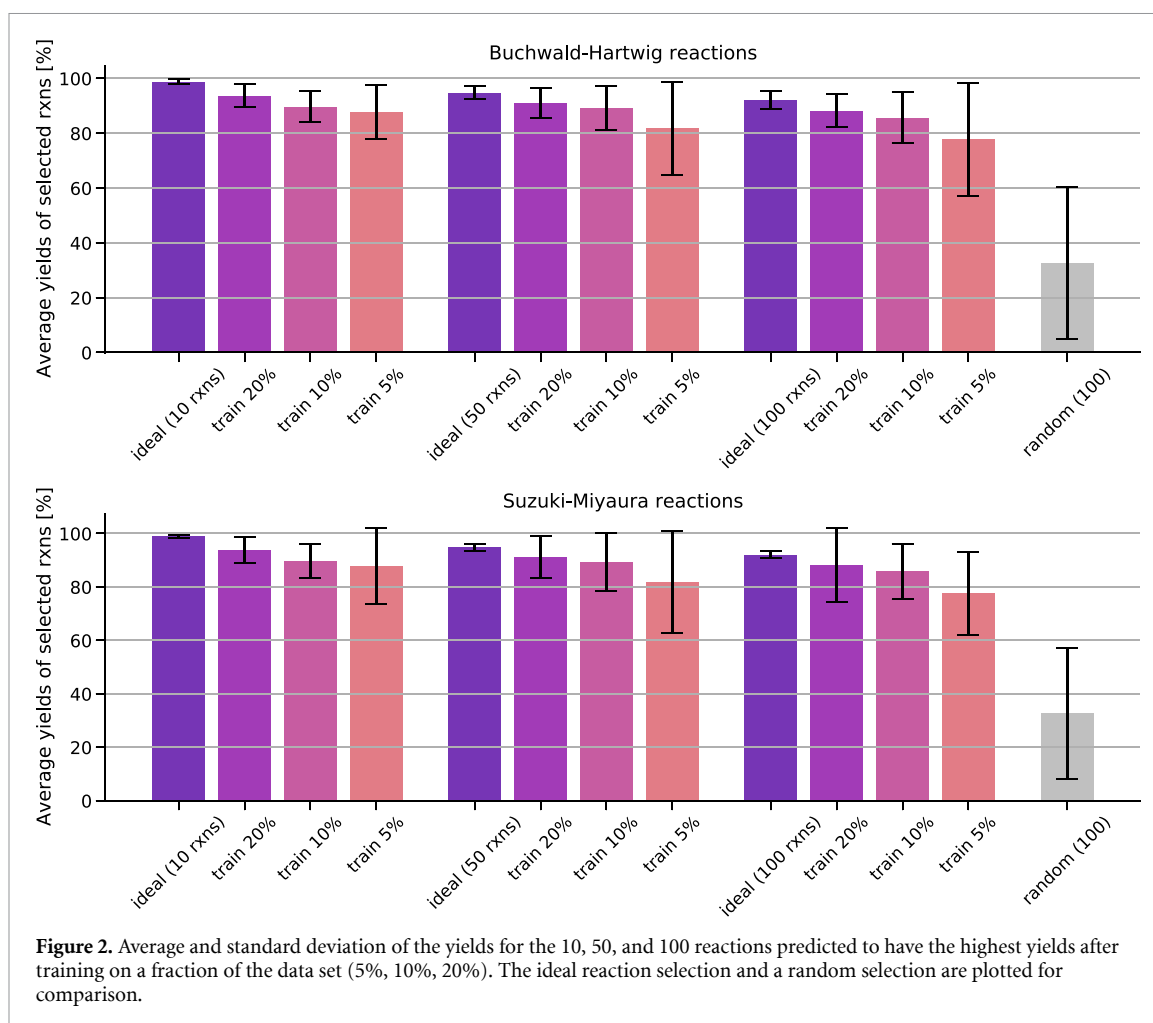


Figure 2. Average and standard deviation of the yields for the 10, 50, and 100 reactions predicted to have the highest yields after training on a fraction of the data set (5%, 10%, 20%). The ideal reaction selection and a random selection are plotted for comparison.

4. Patent yield predictions

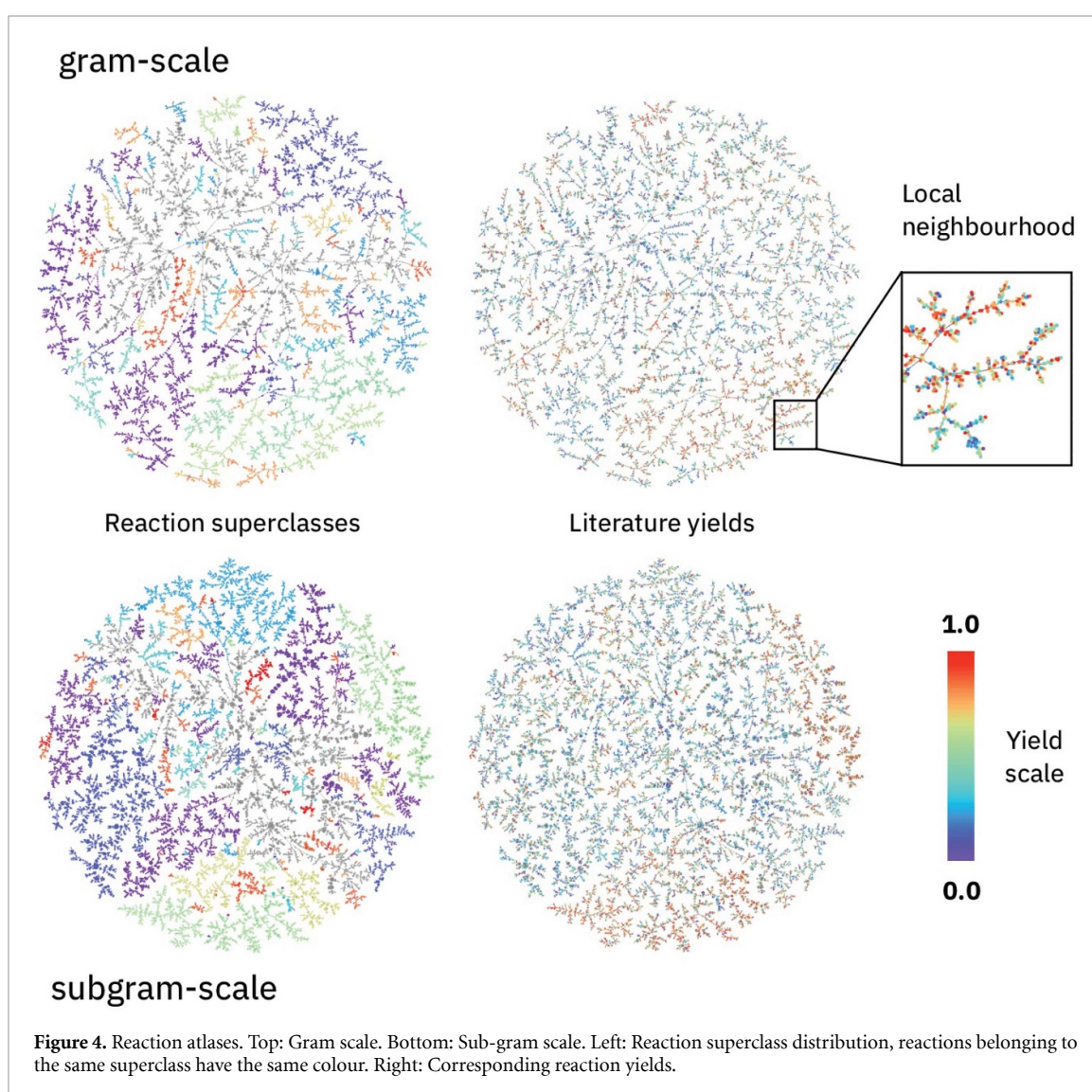
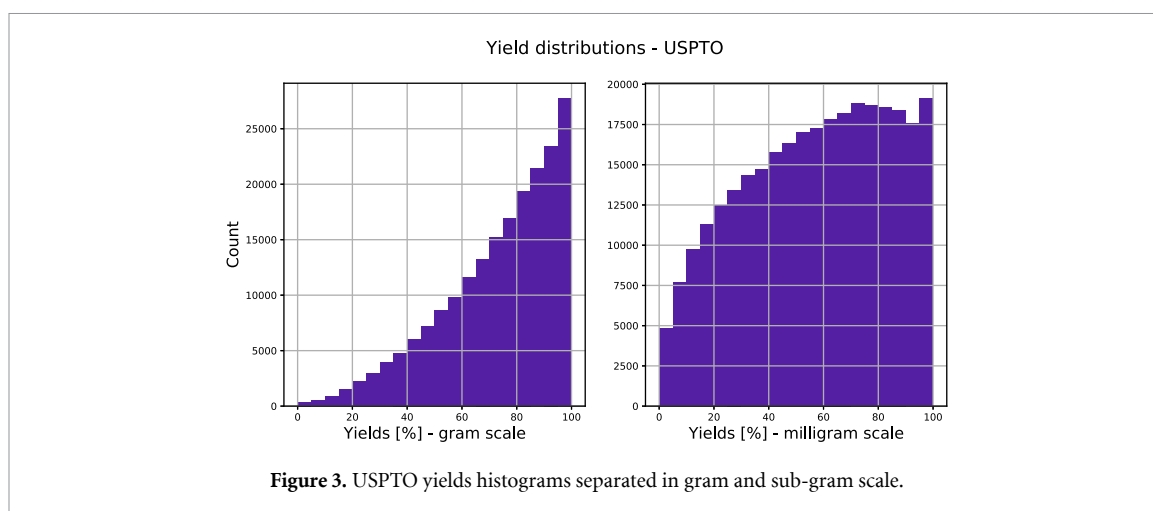
In this section, we analyse USPTO data set [22, 23] yields. We started from the same set as in our previous work [28], keeping only reactions for which yields and product mass were reported. In contrast to HTE, where reactions are typically performed in sub-gram scale, the patent data contains reactions spanning a wider range, from grams to sub-grams scales.

4.1. Gram versus sub-gram scale

When investigating the yields for different mass scales, we observed that gram and sub-gram scales had statistically different yield distributions, as shown in figure 3. Tables S1–S3 show an additional analysis of the two scales. One reason could be that the reaction sub-gram scale reactions are generally less optimised than gram-scale. In sub-gram scale, the primary goal is to show that the desired product is present. To be able to synthesise a specific compound on a larger scale, reactions are optimised and predominantly high-yield reactions are employed. Therefore, we split the USPTO reactions into two data sets according to the product mass. If for the same canonical reaction SMILES multiple yields were reported in the same mass scale, we took the average of those yields.

We performed various experiments, as summarised in table 3. The R^2 scores for the randomly train-test splits with 0.117 for gram scale and 0.195 low. As expected, the tasks become even more difficult when the time split is used. In our experiment, we took all reactions first published in 2012 and before as a training/validation set and the reactions published after 2012 as a test set. To show that the model was still able to learn, we performed a sanity check by randomising the yields across the training reactions. The resulting performance on the test set was a R^2 score of 0.

Unfortunately, the yields from the USPTO data set could not be accurately predicted. To better understand why, we further inspected the USPTO reaction yields with a visual analysis using reaction atlases built using TMAP [29], faerun [30] and our reaction fingerprints [17]. Figure 4 reveals that globally reaction classes tend to have similar yields. However, if a local neighbourhood is analysed the nearest neighbours



often have extremely diverse reaction yields. Those diverse yields make it challenging for the model to learn anything but yield averages for similar reactions and hence, explain the low performance on the patent reactions. This analysis opens up relevant questions on the quality of the reported information (relative to the mass scale) and its extraction accuracy from text, which could severely hamper the development of reaction yield predictive models. The need of cleaned and consistent reaction yields data set is even more important than for other reaction prediction tasks.

Table 3. Summary of the R^2 scores on the different USPTO reaction sets.

Scale	Gram	Sub-gram
Random split	0.117	0.195
Time split	0.095	0.142
Random split (smoothed)	0.277	0.388
Randomized yields	0.0	0.0

In table 3, the ‘random split (smoothed)’ row shows an experiment inspired from the observations above. As some of the yields values are probably incorrect in the data set, we smoothed the yields by computing the average of the three nearest neighbour yields plus twice the own yield of the reaction. The nearest neighbours were estimated using the *rxnfp ft* [17] and *faiss* [31]. Figure S25 shows the yield distributions after smoothing. On the smoothed data sets, the performance of our models more than triples in the gram scale and doubles on the sub-gram scale, achieving R^2 scores of 0.277 and 0.388, respectively. The removal of noisy reactions [32] or reaction data augmentation techniques [33] could potentially lead to further improvements.

5. Conclusion

In this work, we combined a reaction SMILES encoder with a reaction regression task to design a reaction yield predictive model. We analysed two HTE reaction data sets, showing excellent results. On the Buchwald–Hartwig reaction data set, our models outperform previous work on random splits and perform similar to models trained on chemical descriptors computed with DFT on test sets where specific additives were held out from the training set. Compared to random forest models, the feature importance can not directly be obtained. Future work could (visually) investigate the attention weights to find out what tokens and molecules contribute the most to the predictions [34, 35].

We analysed the yields in the public patent data and show that the distribution of reported yields strongly differs depending on the reaction scale. Because of the intrinsic lack of consistency and quality in the patent data, our proposed method fails to predict patent reaction yields accurately. While we cannot rule out the existence of any other architecture potentially performing better than the one presented in this manuscript, we raise the need for a more consistent and better quality public data set for the development of reaction yields prediction models. The suspicion that the patent data yields are inconsistently reported is substantiated by the large variability of methods used to purify and report yields by the different reaction mass scales and the different optimisation in each reported reaction. Our reaction atlases [17, 29, 30] reveal globally higher yielding reaction classes. However, nearest neighbours often have significantly scattered yields. We show that better results can be achieved by smoothing the patent data yields using the nearest neighbours.

Our approach to yield predictions can be extended to any reaction regression task, for example, for predicting reaction activation energies [36–38], and is expected to have a broad impact in the field of organic chemistry.

Code and data statement

The code and data are available on https://rxn4chemistry.github.io/rxn_yields/.

Acknowledgment

We acknowledge the RXN for Chemistry team for insightful discussion.

ORCID iDs

Philippe Schwaller  <https://orcid.org/0000-0003-3046-6576>

Alain C Vaucher  <https://orcid.org/0000-0001-7554-0288>

Teodoro Laino  <https://orcid.org/0000-0001-8717-0456>

Jean-Louis Reymond  <https://orcid.org/0000-0003-2724-2942>

References

- [1] Schwaller P, Hoover B, Reymond J-L, Strobel H and Laino T 2020 Unsupervised Attention-Guided Atom-Mapping ChemRxiv preprint (<https://doi.org/10.26434/chemrxiv.12298559.v1>)

- [2] Segler M H, Preuss M and Waller M P 2018 Planning chemical syntheses with deep neural networks and symbolic AI *Nature* **555** 604–10
- [3] Coley C W *et al* 2019 A robotic platform for flow synthesis of organic compounds informed by AI planning *Science* **365** eaax1566
- [4] Schwaller P *et al* 2020 Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy *Chem. Sci.* **11** 3316–25
- [5] Genheden S *et al* 2020 AiZynthFinder: a fast, robust and flexible open-source software for retrosynthetic planning *J. Cheminform.* **12** 70
- [6] Schwaller P *et al* 2019 Molecular transformer: a model for uncertainty-calibrated chemical reaction prediction *ACS Cent. Sci.* **5** 1572–83
- [7] Kite S, Hattori T and Murakami Y 1994 Estimation of catalytic performance by neural network—product distribution in oxidative dehydrogenation of ethylbenzene *Appl. Catal. A* **114** L173–78
- [8] Raccuglia P *et al* 2016 Machine-learning-assisted materials discovery using failed experiments *Nature* **533** 73–6
- [9] Ahneman D T, Estrada J G, Lin S, Dreher S D and Doyle A G 2018 Predicting reaction performance in C–N cross-coupling using machine learning *Science* **360** 186–90
- [10] Chuang K V and Keiser M J 2018 Comment on “Predicting reaction performance in C–N cross-coupling using machine learning” *Science* **362** 6416
- [11] Sandfort F, Strieth-Kalthoff F, Kühnemund M, Beecks C and Glorius F 2020 A structure-based platform for predicting chemical reactivity *Chem.* **6** 1379–90
- [12] Granda J M, Donina L, Dragone V, Long D-L and Cronin L 2018 Controlling an organic synthesis robot with machine learning to search for new reactivity *Nature* **559** 377–81
- [13] Fu Z *et al* 2020 Optimizing chemical reaction conditions using deep learning: a case study for the Suzuki–Miyaura cross-coupling reaction *Org. Chem. Front.* **7** 2269–77
- [14] Eyke N S, Green W H and Jensen K F 2020 Iterative experimental design based on active machine learning reduces the experimental burden associated with reaction screening *React. Chem. Eng.* **5** 1963–72
- [15] Skoraczynski G *et al* 2017 Predicting the outcomes of organic reactions via machine learning: are current descriptors sufficient? *Sci. Rep.* **7** 3582
- [16] Schwaller P, Gaudin T, Lanyi D, Bekas C and Laino T 2018 “Found in translation”: predicting outcomes of complex organic chemistry reactions using neural sequence-to-sequence models *Chem. Sci.* **9** 6091–8
- [17] Schwaller P *et al* 2021 Mapping the space of chemical reactions using attention-based neural networks *Nat. Mach. Intell.* **3** 144–52
- [18] Devlin J, Chang M-W, Lee K and Toutanova K 2019 BERT: pre-training of deep bidirectional transformers for language understanding *Proc. of the 2019 Conf. of the North {A}merican Chapter of the Association for Computational Linguistics: Human Language Technologies* (Stroudsburg, PA: Association for Computational Linguistics) 4171–86
- [19] Vaswani A *et al* 2017 Attention is all you need *Advances in Neural Information Processing Systems 30* (Red Hook, NY: Curran Associates, Inc.) 5998–6008 (<https://papers.nips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>)
- [20] Weininger D 1988 SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules *J. Chem. Inf. Model.* **28** 31–6
- [21] Perera D *et al* 2018 A platform for automated nanomole-scale reaction screening and micromole-scale synthesis in flow *Science* **359** 429–34
- [22] Lowe D M 2012 Extraction of chemical structures and reactions from the literature *PhD Thesis* University of Cambridge (<https://doi.org/10.17863/CAM.16293>)
- [23] Lowe D 2017 Chemical reactions from US patents (1976–Sep2016) (<https://doi.org/10.6084/m9.figshare.5104873.v1>)
- [24] Simpletransformers (available at: <https://simpletransformers.ai>) (Accessed: 2 July 2020)
- [25] Wolf T *et al* 2020 Transformers: State-of-the-art natural language processing *Proc. 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations* (Stroudsburg, PA: Association for Computational Linguistics) 38–45
- [26] Paszke A *et al* 2019 PyTorch: an imperative style, high-performance deep learning library *Proc. of Advances in Neural Information Processing Systems 32* (Red Hook, NY: Curran Associates, Inc.) 8026–37 (<https://papers.nips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html>)
- [27] Landrum G *et al* 2019 rdkit/rdkit: 2019_03_4 (q1 2019) release (<https://doi.org/10.5281/zenodo.3366468>)
- [28] Pesciullesi G, Schwaller P, Laino T and Reymond J-L 2020 Transfer learning enables the molecular transformer to predict regio- and stereoselective reactions on carbohydrates *Nat. Commun.* **11** 1–8
- [29] Probst D and Reymond J-L 2020 Visualization of very large high-dimensional data sets as minimum spanning trees *J. Cheminform.* **12** 1–13
- [30] Probst D and Reymond J-L 2017 Fun: a framework for interactive visualizations of large, high-dimensional datasets on the web *Bioinformatics* **34** 1433–5
- [31] Johnson J, Douze M and Jégou H 2019 Billion-scale similarity search with GPUs *IEEE Trans. Big Data* (<https://doi.org/10.1109/TBDDATA.2019.2921572>)
- [32] Toniato A, Schwaller P, Cardinale A, Geluykens J and Laino T 2020 Unassisted noise-reduction of chemical reactions data sets ChemRxiv preprint (<https://doi.org/10.26434/chemrxiv.12395120.v1>)
- [33] Tetko I V, Karpov P, Van Deursen R and Godin G 2020 State-of-the-art augmented NLP transformer models for direct and single-step retrosynthesis *Nat. Commun.* **11** 5575
- [34] Hoover B, Strobel H and Gehrmann S 2019 Exbert: a visual analysis tool to explore learned representations in transformers models *Proc. 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations* (Stroudsburg, PA: Association for Computational Linguistics) 187–96
- [35] Vig J and Belinkov Y 2019 Analyzing the structure of attention in a transformer language model *Proc. of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP* (Stroudsburg, PA: Association for Computational Linguistics) 63–76
- [36] Grambow C A, Pattanaik L and Green W H 2020 Reactants, products and transition states of elementary chemical reactions based on quantum chemistry *Sci. Data* **7** 1–8
- [37] von Rudorff G F, Heinen S, Bragato M and von Lilienfeld A 2020 Thousands of reactants and transition states for competing E2 and SN2 reactions *Mach. Learn.: Sci. Technol.* **1** 045026
- [38] Jorner K, Brinck T, Norrby P-O and Buttar D 2020 Machine learning meets mechanistic modelling for accurate prediction of experimental activation energies *Chem. Sci.* **12** 1163–75